

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

Install Pyspark

```
In [1]: !pip install pyspark
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting pyspark
  Downloading pyspark-3.2.1.tar.gz (281.4 MB)
    |████████████████████████████████████████| 281.4 MB 33 kB/s
Collecting py4j==0.10.9.3
  Downloading py4j-0.10.9.3-py2.py3-none-any.whl (198 kB)
    |████████████████████████████████████████| 198 kB 41.7 MB/s
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.2.1-py2.py3-none-any.whl size=281853642 sha256=f2686fb889f8e1813d68af554ef1a8228ef3daf58a57f3c972b749190245ab42
  Stored in directory: /root/.cache/pip/wheels/9f/f5/07/7cd8017084dce4e93e84e92efd1e1d5334db05f2e83bcef74f
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully installed py4j-0.10.9.3 pyspark-3.2.1
```

Create Session

The entry point to programming Spark with the Dataset

```
In [2]: from pyspark.sql import SparkSession
```

```
In [3]: session = SparkSession.builder.appName("sql").master("local").getOrCreate()
```

```
In [5]: data = session.read.csv("churn.csv", header = True, inferSchema = True)
```

```
In [7]: data.show(5)
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+
--+-+-----+-----+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+
|customerID|gender|SeniorCitizen|Partner|Dependents|tenure|CallService|MultipleConnections|InternetConnection|OnlineSecurity|OnlineBackup|DeviceProtectionService|TechnicalHelp|OnlineTV|OnlineMovies|Agreement|BillingMethod|PaymentMethod|MonthlyServiceCharges|TotalAmount|Churn|
+-----+-----+-----+-----+-----+-----+-----+-----+
--+-+-----+-----+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+
|2907-ILJBN|Female|0.0|Yes|Yes|11.0|Yes|
No|No|No internet service|No internet service|No internet service|No internet service|No internet service|No internet service|One year|No|
```



```
In [22]: from pyspark.sql.functions import count, when, col
fourth = data.filter(data['MonthlyServiceCharges'] > 100).groupBy("gender").agg(count(when(
+-----+-----+-----+-----+
+-----+
|gender|count(CASE WHEN (Dependents = Yes) THEN True END)|count(CASE WHEN (Dependents =
No) THEN True END)|
+-----+-----+-----+-----+
+-----+
|Female|                                     183|
|         438|
|   Male|                                     182|
|         389|
+-----+-----+-----+-----+
+-----+
```

Q5 How many number of customers have churned and not churned. Consider only female customers who have no dependents and has done call service and has preferred electronic check method.

```
In [33]: from pyspark.sql.functions import count
fifth = data.filter((data['gender']=='Female') & (data['Dependents']=='No') & (data['Cal
+-----+-----+-----+-----+
+-----+
|count(CASE WHEN (Churn = Yes) THEN True END)|count(CASE WHEN (Churn = No) THEN True EN
D)|
+-----+-----+-----+-----+
+-----+
|                                     718|                                     56
6|
+-----+-----+-----+-----+
+-----+
```

Q6 How many male and female customers have no dependents and have multiple connections. Consider the customers who have call service and has preferred either electronic check method or mailed check method

```
In [24]: from pyspark.sql.functions import count
sixth = data.filter((data['CallService'] == 'Yes') & ((data['PaymentMethod'] == 'Electro
+-----+-----+-----+-----+
+-----+
|gender|count(CASE WHEN (Dependents = No) THEN True END)|count(CASE WHEN (MultipleConnec
tions = Yes) THEN True END)|
+-----+-----+-----+-----+
+-----+
|Female|                                     2068|
|         1252|
|   Male|                                     2071|
|         1233|
+-----+-----+-----+-----+
+-----+
```

Q8 What is the maximum monthly service charges of customers who have done payment by electronic check method? Consider only those customers who have agreement for on year or two years only.

```
In [25]: from pyspark.sql.functions import max
eight = data.filter((data['PaymentMethod'] == 'Electronic check') & ((data['Agreement']
+-----+-----+
|max(MonthlyServiceCharges)|
+-----+
|                                     118.65|
```

+-----+

Q9 What is the minimum total amount of male and female customers having one year or two year agreement type. Consider only those customers who have no internet connection, no online security no online backup and no device protection service.

```
In [27]: from pyspark.sql.functions import min
ninth = data.filter(((data['Agreement'] == 'One year') | (data['Agreement'] == 'Two year

+-----+-----+
|gender|min(TotalAmount)|
+-----+-----+
|Female|      69.21978888|
|  Male|      59.02463086|
+-----+-----+
```

```
In [ ]:
```